

PCT

WORLD INTELLECTUAL PROPERTY ORGANIZATION
International Bureau



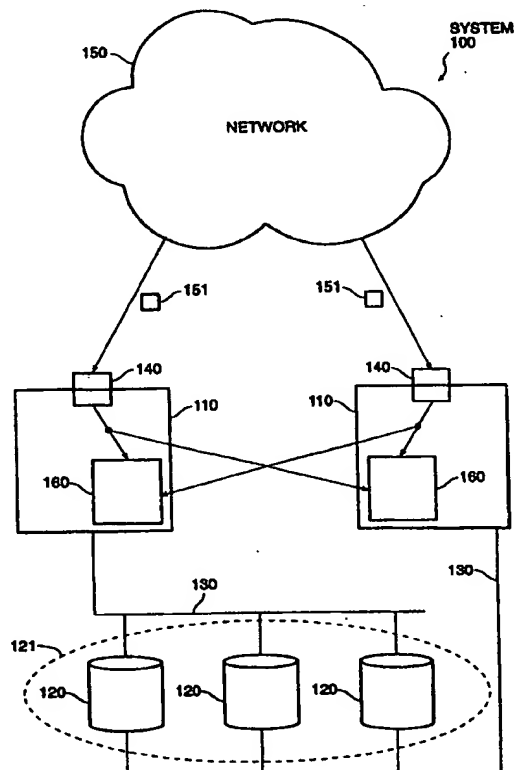
INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(51) International Patent Classification 6 : G06F 11/14	A1	(11) International Publication Number: WO 99/46680 (43) International Publication Date: 16 September 1999 (16.09.99)
(21) International Application Number: PCT/US99/05071 (22) International Filing Date: 8 March 1999 (08.03.99) (30) Priority Data: 09/037,652 10 March 1998 (10.03.98) US (71) Applicant: NETWORK APPLIANCE, INC. [US/US]; 2770 San Tomas Expressway, Santa Clara, CA 95051 (US). (72) Inventor: KLEINMAN, Steven; 157 El Monte Court, Los Altos, CA 94022 (US). (74) Agent: LAW OFFICES OF STEVEN A. SWERNOFSKY; P.O. Box 390013, Mountain View, CA 94039-0013 (US).		(81) Designated States: CA, CN, JP, KR, European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE). Published <i>With international search report. Before the expiration of the time limit for amending the claims and to be republished in the event of the receipt of amendments.</i>

(54) Title: **HIGHLY AVAILABLE FILE SERVERS**

(57) Abstract

The invention provides a storage system that is highly available even in the face of component failures in the storage system, and a method for operating that storage system. A first and a second file server each includes a file server request log for storing incoming file server requests. Both the first and second file servers have access to a common set of mass storage elements. Each incoming file server request is copied to both the first and second file servers; the first file server processes the file server request while the second file server maintains a copy in its file server request log. Each file server operates using a file system that maintains consistent state after each file server request. On failover, the second file server can perform those file server requests in its file server request log since the most recent consistent state. There is no single point of failure that prevents access to any individual mass storage element.



BEST AVAILABLE COPY

FOR THE PURPOSES OF INFORMATION ONLY

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

AL	Albania	ES	Spain	LS	Lesotho	SI	Slovenia
AM	Armenia	FI	Finland	LT	Lithuania	SK	Slovakia
AT	Austria	FR	France	LU	Luxembourg	SN	Senegal
AU	Australia	GA	Gabon	LV	Latvia	SZ	Swaziland
AZ	Azerbaijan	GB	United Kingdom	MC	Monaco	TD	Chad
BA	Bosnia and Herzegovina	GE	Georgia	MD	Republic of Moldova	TG	Togo
BB	Barbados	GH	Ghana	MG	Madagascar	TJ	Tajikistan
BE	Belgium	GN	Guinea	MK	The former Yugoslav Republic of Macedonia	TM	Turkmenistan
BF	Burkina Faso	GR	Greece			TR	Turkey
BG	Bulgaria	HU	Hungary	ML	Mali	TT	Trinidad and Tobago
BJ	Benin	IE	Ireland	MN	Mongolia	UA	Ukraine
BR	Brazil	IL	Israel	MR	Mauritania	UG	Uganda
BY	Belarus	IS	Iceland	MW	Malawi	US	United States of America
CA	Canada	IT	Italy	MX	Mexico	UZ	Uzbekistan
CF	Central African Republic	JP	Japan	NE	Niger	VN	Viet Nam
CG	Congo	KE	Kenya	NL	Netherlands	YU	Yugoslavia
CH	Switzerland	KG	Kyrgyzstan	NO	Norway	ZW	Zimbabwe
CI	Côte d'Ivoire	KP	Democratic People's Republic of Korea	NZ	New Zealand		
CM	Cameroon		Republic of Korea	PL	Poland		
CN	China	KR	Republic of Korea	PT	Portugal		
CU	Cuba	KZ	Kazakstan	RO	Romania		
CZ	Czech Republic	LC	Saint Lucia	RU	Russian Federation		
DE	Germany	LI	Liechtenstein	SD	Sudan		
DK	Denmark	LK	Sri Lanka	SE	Sweden		
EE	Estonia	LR	Liberia	SG	Singapore		

Title of the Invention

Highly Available File Servers

Background of the Invention

1. Field of the Invention

The invention relates to storage systems.

2. Related Art

Computer storage systems are used to record and retrieve data. In some computer systems, storage systems communicate with a set of client devices, and provide services for recording and retrieving data to those client devices. Because data storage is important to many applications, it is desirable for the services and data provided by the storage system to be available for service to the greatest degree possible. It is therefore desirable to provide storage systems that can remain available for service even in the face of component failures in the storage system.

One known technique for provide storage systems that can remain available for service is to provide a plurality of redundant storage elements, with the property that when a first storage element fails, a second storage element is available to provide the services and the data otherwise provided by the first. Transfer of the function of providing services from the first to the second storage element is called "failover." The second storage element maintains a copy of the data maintained by the first, so that failover can proceed without substantial interruption.

A first known technique for achieving failover is to cause the second storage element to copy all the operations of the first. Thus, each storage operation completed by the first storage element is also completed by the second. This first known technique is subject to drawbacks: (1) It uses a substantial amount of processing power at the second storage element duplicating efforts of the first, most of which is wasted. (2) It slows the first storage element in confirming completion of operations, because the first storage element waits for the second to also complete the same operations.

1 A second known technique for achieving failover is to identify a sequence of
2 checkpoints at which the first storage element is at a consistent and known state. On failover,
3 the second storage element can continue operation from the most recent checkpoint. For exam-
4 ple, the NFS (Network File System) protocol requires all write operations to be stored to disk
5 before they are confirmed, so that confirmation of a write operation indicates a stable file system
6 configuration. This second known technique is subject to drawbacks: (1) It slows the first stor-
7 age element in performing write operations, because the first storage element waits for write op-
8 erations to be completely stored to disk. (2) It slows recovery on failover, because the second
9 storage element addresses any inconsistencies left by failure of the first between identified
10 checkpoints.

11

12 Accordingly, it would be advantageous to provide a storage system, and a method
13 for operating a storage system, that efficiently uses all storage system elements, quickly com-
14 pletes and confirms operations, and quickly recovers from failure of any storage element. This
15 advantage is achieved in an embodiment of the invention in which the storage system imple-
16 ments frequent and rapid checkpoints, and in which the storage system rapidly distributes dupli-
17 cate commands for those operations between checkpoints among its storage elements.

18

19

Summary of the Invention

20

21 The invention provides a storage system that is highly available even in the face
22 of component failures in the storage system, and a method for operating that storage system. A
23 first and a second file server each includes a file server request log for storing incoming file
24 server requests. Both the first and second file servers have access to a common set of mass stor-
25 age elements. Each incoming file server request is copied to both the first and second file serv-
26 ers; the first file server processes the file server request while the second file server maintains a
27 copy in its file server request log. Each file server operates using a file system that maintains
28 consistent state after each file server request. On failover, the second file server can perform
29 those file server requests in its file server request log since the most recent consistent state.

30

31 In a second aspect of the invention, a file server system provides mirroring of one
32 or more mass storage elements. Each incoming file server request is copied to both the first file
33 server and the second file server. The first file server performs the file server requests to modify
34 a primary set of mass storage elements, and also performs the same file server requests to mod-
35 ify a mirror set of mass storage elements. The mirror mass storage elements are disposed physi-

1 cally separately from the primary mass storage elements, such as at another site, and provide a
2 resource in the event the entire primary set of mass storage elements is to be recovered.

3 4 Brief Description of the Drawings

5
6 Figure 1 shows a block diagram of a highly available file server system.

7
8 Figure 2 shows a block diagram of a file server in the file server system.

9
10 Figure 3 shows a process flow diagram of operation of the file server system.

11 12 Detailed Description of the Preferred Embodiment

13
14 In the following description, a preferred embodiment of the invention is described
15 with regard to preferred process steps and data structures. However, those skilled in the art
16 would recognize, after perusal of this application, that embodiments of the invention may be im-
17 plemented using one or more general purpose processors (or special purpose processors adapted
18 to the particular process steps and data structures) operating under program control, and that im-
19 plementation of the preferred process steps and data structures described herein using such
20 equipment would not require undue experimentation or further invention.

21 22 *File Server Pair and Failover Operation*

23
24 Figure 1 shows a block diagram of a highly available file server system.

25
26 A file server system 100 includes a pair of file servers 110, both coupled to a
27 common set of mass storage devices 120. A first one of the file servers 110 is coupled to a first
28 I/O bus 130 for controlling a first selected subset of the mass storage devices 120. Similarly, a
29 second one of the file servers 110 is coupled to a second I/O bus 130 for controlling a second
30 selected subset of the mass storage devices 120.

31
32 Although both file servers 110 are coupled to all of the common mass storage de-
33 vices 120, only one file server 110 operates to control any one mass storage device 120 at any
34 designated time. Thus, even though the mass storage devices 120 are each controllable by only

1 one file server 110 at a time, each of the mass storage devices 120 remains available even if one
2 of its two associated file servers 110 fails.

3
4 In a preferred embodiment, the file server system 100 includes a pair of such file
5 servers 110; however, in alternative embodiments, more than two such file servers 110 may be
6 included in a single file server system 100.

7
8 In a preferred embodiment, the first I/O bus 130 and the second I/O bus 130 each
9 include a mezzanine bus such as the PCI bus architecture.

10
11 In a preferred embodiment, the mass storage devices 120 include magnetic disk
12 drives, optical disk drives, or magneto-optical disk drives. In alternative embodiments, however,
13 other storage systems may be used, such as bubble memory, flash memory, or systems using
14 other storage technologies. Components of the mass storage devices 120 are referred to as
15 "disks," even though those components may comprise other forms or shapes.

16
17 Each mass storage device 120 can include a single disk or a plurality of disks. In
18 a preferred embodiment, each mass storage device 120 includes a plurality of disks and is dis-
19 posed and operated as a RAID storage system.

20
21 In a preferred embodiment, the first file server 110 is coupled to the second file
22 server 110 using a common interconnect. The common interconnect provides a remote memory
23 access capability for each file server 110, so that data can be stored at each file server 110 from a
24 remote location. In a preferred embodiment, the common interconnect includes a Tandem
25 "ServerNet" interconnect. The common interconnect is coupled to each file server 110 using a
26 device controller coupled to an I/O bus for each file server 110.

27
28 The first file server 110 is coupled to a first network interface 140, which is dis-
29 posed to receive file server requests 151 from a network 150. Similarly, the second file server
30 110 is coupled to a second network interface 140, which is also disposed to receive file server
31 requests 151 from the network 150.

32
33 The first file server 110 includes a first server request memory 160, which re-
34 ceives the file server requests 151 and records them. In the event the first file server 110 recov-
35 ers from a power failure or other service disruption, the outstanding file server requests 151 in

1 the first server request memory 160 are re-performed to incorporate them into a next consistent
2 state of the file system maintained by the first file server 110.

3

4 Similarly, the second file server 110 includes a second server request memory
5 160, which receives the file server requests 151 and records them. In the event the second file
6 server 110 recovers from a power failure or other service disruption, the outstanding file server
7 requests 151 in the second server request memory 160 are re-performed to incorporate them into
8 a next consistent state of the file system maintained by the second file server 110.

9

10 When the first file server 110 receives a file server request 151 from the network
11 150, that file server request 151 is copied into the first server request memory 160. The file
12 server request 151 is also copied into the second server request memory 160 using remote mem-
13 ory access over the common interconnect. Similarly, when the second file server 110 receives a
14 file server request 151 from the network 150, that file server request 151 is copied into the sec-
15 ond server request memory 160. The file server request 151 is also copied into the first server
16 request memory 160 using remote memory access over the common interconnect. Using remote
17 memory access is relatively quicker and has less communication overhead than using a net-
18 working protocol.

19

20 In the event that either file server 110 fails, the other file server 110 can continue
21 processing using the file server requests 151 stored in its own server request memory 160.

22

23 In a preferred embodiment, each server request memory 160 includes a nonvola-
24 tile memory, so those file server requests stored in either server request memory 160 are not lost
25 due to power failures or other service interruptions.

26

27 The responding file server 110 processes the file server request 151 and possibly
28 modifies stored files on one of the mass storage devices 120. The non-responding file server
29 110, partner to the responding file server 110, maintains the file server request 151 stored in its
30 server request memory 160 to prepare for the possibility that the responding file server 110
31 might fail. In the event the responding file server 110 fails, the non-responding file server 110
32 processes the file server request 151 as part of a failover technique.

33

1 In a preferred embodiment, each file server 110 controls its associated mass stor-
2 age devices 120 so as to form a redundant array, such as a RAID storage system, using inven-
3 tions described in the following patent applications:

4
5 o Application Serial No. 08/471,218, filed June 5, 1995, in the name of inventors David
6 Hitz et al., titled "A Method for Providing Parity in a Raid Sub-System Using Non-
7 Volatile Memory", attorney docket number NET-004;

8
9 o Application Serial No. 08/454,921, filed May 31, 1995, in the name of inventors David
10 Hitz et al., titled "Write Anywhere File-System Layout", attorney docket number NET-
11 005;

12
13 o Application Serial No. 08/464,591, filed May 31, 1995, in the name of inventors David
14 Hitz et al., titled "Method for Allocating Files in a File System Integrated with a Raid
15 Disk Sub-System", attorney docket number NET-006.

16
17 Each of these applications is hereby incorporated by reference as if fully set forth
18 herein. They are collectively referred to as the "WAFL Disclosures."

19
20 As part of the techniques shown in the WAFL Disclosures, each file server 110
21 controls its associated mass storage devices 120 in response to file server requests 151 in an
22 atomic manner. The final action for any file server request 151 is to incorporate the most recent
23 consistent state into the file system 121. Thus, file system 121 is in an internally consistent state
24 after completion of each file server request 151. Thus, a file system 121 defined over the mass
25 storage devices 120 will be found in an internally consistent state, regardless of which file server
26 110 controls those mass storage devices 120. Exceptions to the internally consistent state will
27 only include a few of the most recent file server requests 151, which will still be stored in the
28 server request memory 160 for both file servers 110. Those most recent file server requests 151
29 can be incorporated into a consistent state by performing them with regard to the most recent
30 consistent state.

31
32 For any file server request 151, in the event the file server 110 normally re-
33 sponding to that file server request 151 fails, the other file server 110 will recognize the failure
34 and perform a failover method to take control of mass storage devices 120 previously assigned
35 to the failing file server 110. The failover file server 110 will find those mass storage devices

1 120 with their file system 121 in an internally consistent state, but with the few most recent file
2 server requests 151 as yet unperformed. The failover file server 110 will have copies of these
3 most recent file server requests 151 in its server request memory 160, and will perform these file
4 server requests 151 in response to those copies.

5

6 *File Server Node*

7

8 Figure 2 shows a block diagram of a file server in the file server system.

9

10 Each file server 110 includes at least one processor 111, a program and data
11 memory 112, the server request memory 160 (including a nonvolatile RAM), a network interface
12 element 114, and a disk interface element 115. These elements are interconnected using a bus
13 117 or other known system architecture for communication among processors, memory, and pe-
14 ripherals.

15

16 In a preferred embodiment, the network interface element 114 includes a known
17 network interface for operating with the network 150. For example, the network interface ele-
18 ment 114 can include an interface for operating with the FDDI interface standard or the
19 100BaseT interface standard.

20

21 After failover, the file server 110 responds to file server requests directed to either
22 itself or its (failed) partner file server 110. Each file server 110 is therefore capable of assuming
23 an additional network identity on failover, one for itself and one for its failed partner file server
24 110. In a preferred embodiment, the network interface element 114 for each file server 110 in-
25 cludes a network adapter capable of responding to two separate addresses upon instruction by
26 the file server 110. In an alternative embodiment, each file server 110 may have two such net-
27 work adapters.

28

29 In a preferred embodiment, the disk interface element 115 includes a known disk
30 interface for operating with magnetic, optical, or magneto-optical disks, that has two independ-
31 ent ports with each port coupled to a separate file server 110, such as the FC-AL interface. This
32 helps prevent failure of one file server 110 from affecting low-level operation of the other file
33 server 110.

34

1 In a preferred embodiment, the bus 117 includes at least a memory bus 171 and
2 the mezzanine bus 130. The memory bus 171 couples the processor 111 and the program and
3 data memory 112. The mezzanine bus 130 couples the network interface element 114 and the
4 disk interface element 115. The memory bus 171 is coupled to the mezzanine bus 130 using an
5 I/O controller 173 or other known bus adapter technique.

6
7 In a preferred embodiment, each disk in the mass storage 120 is statically as-
8 signed to either the first file server 110 or the second file server 110, responsive to whether the
9 disk is wired for primary control by either the first file server 110 or the second file server 110.
10 Each disk has two control ports A and B; the file server 110 wired to port A has primary control
11 of that disk, while the other file server 110 only has control of that disk when the other file
12 server 110 has failed.

13 14 *Operation Process Flow*

15
16 Figure 3 shows a process flow diagram of operation of the file server system.

17
18 A method 300 is performed by the components of the file server 100, and in-
19 cludes a set of flow points and process steps as described herein.

20
21 At a flow point 310, a device coupled to the network 150 desires to make a file
22 system request 151.

23
24 At a step 311, the device transmits the file system request 151 to the network 150.

25
26 At a step 312, the network 150 transmits the file server request 151 to the file
27 server 110.

28
29 At a step 313, a first file server 110 at the file server system 100 receives the file
30 server request 151. The first file server 110 copies the file server request 151 into the first server
31 request memory 160, and also copies the file server request 151 into the second server request
32 memory 160 using the common interconnect. The target of the copying operation in the second
33 server request memory 160 is to an area reserved for this purpose. The copying operation re-
34 quires no further processing by the second file server 110, and the second file server 110 does
35 not normally process or respond to the file server request 151.

1
2 At a step 314, the first file server 110 responds to the file server request 151.

3
4 At a flow point 320, the file server request has been successfully processed.

5
6 In a second aspect of the invention, the first file server 110 provides mirroring of
7 one or more of its mass storage devices 120.

8
9 As with the first aspect of the invention, each incoming file server request is
10 copied to both the first file server 110 and the second file server 110. The first file server 110
11 performs the file server requests to modify one or more primary mass storage devices 120 under
12 its control. The first file server 110 also performs the file server requests to modify a set of mir-
13 ror mass storage devices 120 under its control, but located distant from the primary mass storage
14 devices 120. Thus, the mirror mass storage devices 120 will be a substantial copy of the primary
15 mass storage devices 120.

16
17 The mirror set of mass storage devices 120 provide a resource in the event the
18 entire primary set of mass storage devices 120 is to be recovered, such as if a disaster befalls the
19 primary set of mass storage devices 120.

20
21 At a flow point 330, the first file server 110 in the file server system 100 fails.

22
23 At a step 331, the second file server 110 in the file server system 100 recognizes
24 the failure of the first file server 110.

25
26 In a preferred embodiment, the second file server 110 performs the step 331 in
27 the following manner:

28
29 o Each file server 110 maintains two disks of its mass storage devices 120 (thus, there are a
30 total of four such disks for two file servers 110) for recording state information about the
31 file server 110. There are two such disks (called "mailbox disks") so that one can be
32 used as primary storage and one can be used as backup storage. If one of the two mail-
33 box disks fails, the file server 110 using that mailbox disk designates another disk as one
34 of its two mailbox disks.

1 o Each file server 110 maintains at least one sector on each mailbox disk, on which the file
2 server 110 periodically writes state information. Each file server 110 also sends its state
3 information to the other file server 110 using the interconnect using remote memory ac-
4 cess. The state information written to the mailbox disks by each file server 110 changes
5 with each update.

6
7 o Each file server 110 periodically reads the state information from at least one of the
8 mailbox disks for the other file server 110. Each file server 110 also receives state in-
9 formation from the other file server 110 using the interconnect using remote memory ac-
10 cess.

11
12 o Each file server 110 recognizes if the other file server 110 has failed by noting that there
13 has been no update to the state information on the mailbox disks for the other file server
14 110.

15
16 In a preferred embodiment, the second file server 110 determines whether failure
17 of the first file server 110 is a hardware error or a software error, and only recognizes failure of
18 the first file server 110 for hardware errors. In alternative embodiments, the second file server
19 110 may recognize failure of the first file server 110 for software errors as well.

20
21 At a step 332, the second file server 110 seizes control of all mass storage devices
22 120 previously assigned to the first file server 110. Due to the nature of the techniques shown in
23 the WAFL Disclosures, the file system 121 defined over those mass storage devices 120 will be
24 in an internally consistent state. All those file server requests 151 marked completed will have
25 been processed and the results incorporated into storage blocks of the mass storage devices 120.

26
27 In normal operation, neither file server 110 places reservations on any of the mass
28 storage devices 120. In the step 332 (only on failover), the second file server 110 seizes control
29 of the mass storage devices 120 previously controlled by the first file server 110, and retains
30 control of those mass storage devices 120 until it is satisfied that the first file server 110 has re-
31 covered.

32
33 When the first file server 110 recovers, it sends a recovery message to the second
34 file server 110. In a preferred embodiment, the second file server 110 relinquishes control of the
35 seized mass storage devices 120 by operator command. However, in alternative embodiments,

1 the second file server 110 may recognize the recovery message from the first file server 110 and
2 relinquish control of the seized mass storage devices 120 in response thereto.

3
4 At a step 333, the second file server 110 notes all file server requests 151 in the
5 area of its server request memory 160 that were copied there by the first file server 110. Those
6 file server requests 151 whose results were already incorporated into storage blocks of the stor-
7 age devices 120 are discarded.

8
9 At a step 334, when the second file server 110 reaches its copy of each file server
10 request 151, the second file server 110 processes the file server request 151 normally.

11
12 At a flow point 340, failover from the first file server 110 to the second file server
13 110 has been successfully handled.

14
15 *Alternative Embodiments*

16
17 Although preferred embodiments are disclosed herein, many variations are possi-
18 ble which remain within the concept, scope, and spirit of the invention, and these variations
19 would become clear to those skilled in the art after perusal of this application.

Claims

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35

1. A file server system including
a first file server including a file server change memory;
a second file server including a file server change memory;
a mass storage element;
said first file server and said second file server being coupled to said mass storage
element;

means for copying a descriptor of a file system change to both said first and second file servers, whereby said first file server processes said file system change while said second file server maintains its copy of said descriptor in its file server change memory; and

means for said second file server to perform a file system change in its file server change memory in response to a service interruption by said first file server.

2. A system as in claim 1, including at least one said mass storage element for each said file server.

3. A system as in claim 1, wherein a first said file server is disposed for processing said file system changes atomically, whereby a second said file server can on failover process exactly those file system changes not already processed by said first file server.

4. A system as in claim 1, wherein a first said file server is disposed to respond identically to service interruptions for itself and for a second said file server.

5. A system as in claim 1, wherein at least one said file server is disposed to delay output to said mass storage element without delaying a response to file system changes.

6. A system as in claim 1, wherein at least one said file server responds to a file system change before committing a result of said file system change to mass storage.

7. A system as in claim 1, wherein
each one of said file servers is coupled to at least a portion of said file server change memory using local memory access; and
each one of said file servers is coupled to at least a portion of said file server change memory using remote memory access.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35

8. A system as in claim 1, wherein said descriptor includes a file server request.

9. A system as in claim 1, wherein said file server change memory includes a disk block.

10. A system as in claim 1, wherein said file server change memory includes a file server request.

11. A system as in claim 1, wherein said file server change memory is disposed to delay output to said mass storage element without delaying a response to file server requests.

12. A system as in claim 1, wherein
said mass storage element includes a file storage system;
each said file server is disposed for leaving said file storage system in an internally consistent state after processing file system changes;
said internally consistent state is associated with a set of completed file system changes;
said set of completed file system changes is identifiable by each said file server.

13. A system as in claim 1, wherein said mass storage element includes a file storage system and each said file server is disposed for leaving said file storage system in an internally consistent state after processing each said file system change.

14. A file server system as in claim 1, wherein
said mass storage element includes a primary mass storage element and a mirror mass storage element; and
said first file server processes said file system change for both said primary mass storage element and said mirror mass storage element.

15. A system as in claim 1, wherein said means for copying includes access to at least one of said first and second file server change memories using a NUMA network.

1 16. A system as in claim 1, wherein said means for copying includes remote
2 memory access to at least one of said first and second file server change memories.
3

4 17. A system as in claim 1, wherein said means for said second file server to
5 perform a file server request in its file server change memory is also operative in response to a
6 service interruption by said second file server.
7

8 18. A file server system including
9 a first file server coupled to a first set of mass storage devices;
10 a second file server coupled to a second set of mass storage devices;
11 a server change memory;
12 said first file server disposed for receiving a file server request and in response
13 thereto copying a descriptor of a file system change into said server change memory; and
14 said first file server disposed for processing said file system change for both said
15 first set of mass storage devices and for at least one said mass storage device in said second set.
16

17 19. A system as in claim 18, wherein
18 said second file server is disposed for receiving a file server request and in re-
19 sponse thereto copying a descriptor of a file system change into said server change memory; and
20 said second file server is disposed for processing said file system change for both
21 said second set of mass storage devices and for at least one said mass storage device in said first
22 set.
23

24 20. A system as in claim 18, wherein said server change memory includes a
25 disk block.
26

27 21. A system as in claim 18, wherein said server change memory includes a
28 file server request.
29

30 22. A system as in claim 18, wherein said server change memory includes a
31 first portion disposed at said first file server and a second portion disposed at said second file
32 server.
33

34 23. A system as in claim 18, wherein

1 said server change memory includes a first portion disposed at said first file
2 server and a second portion disposed at said second file server; and

3 said first file server is disposed for copying said descriptor into both said first
4 portion and said second portion.

5
6 24. A system as in claim 18, wherein
7 said server change memory includes a first portion disposed at said first file
8 server and a second portion disposed at said second file server; and
9 said first file server and said second file server are each disposed for copying said
10 descriptor into both said first portion and said second portion.

11
12 25. A system as in claim 18, wherein said server change memory is disposed
13 to delay output to said mass storage element without delaying a response to file server requests.

14
15 26. A file server system including
16 a plurality of file servers, said plurality of file servers coupled to a mass storage
17 element and at least one file server change memory;
18 each said file server disposed for receiving a file server request and in response
19 thereto copying a descriptor of a file system change into said file server change memory; and
20 each said file server disposed for responding to a service interruption by per-
21 forming a file system change in said file server change memory.

22
23 27. A system as in claim 26, including at least one said mass storage element
24 for each said file server.

25
26 28. A system as in claim 26, including at least one said server change memory
27 for each said file server.

28
29 29. A system as in claim 26, wherein a first said file server is disposed for
30 processing said file system changes atomically, whereby a second said file server can on failover
31 process exactly those file system changes not already processed by said first file server.

32
33 30. A system as in claim 26, wherein a first said file server is disposed to re-
34 spond identically to service interruptions for itself and for a second said file server.

1 31. A system as in claim 26, wherein at least one said file server delays output
2 to said mass storage element without delaying a response to file server requests.

3
4 32. A system as in claim 26, wherein at least one said file server responds to a
5 file system change before committing a result of said file system change to mass storage.

6
7 33. A system as in claim 26, wherein
8 each one of said file servers is coupled to at least a portion of said file server
9 change memory using local memory access; and
10 each one of said file servers is coupled to at least a portion of said file server
11 change memory using remote memory access.

12
13 34. A system as in claim 26, wherein each said file server is disposed for
14 copying said descriptors using a NUMA network.

15
16 35. A system as in claim 26, wherein each said file server is disposed for
17 copying said descriptors using remote memory access.

18
19 36. A system as in claim 26, wherein said file server change memory includes
20 a disk block.

21
22 37. A system as in claim 26, wherein said file server change memory includes
23 a file server request.

24
25 38. A system as in claim 26, wherein said file server change memory is dis-
26 posed to delay output to said mass storage element without delaying a response to file server re-
27 quests.

28
29 39. A system as in claim 26, wherein said mass storage element includes a file
30 storage system and each said file server is disposed for leaving said file storage system in an in-
31 ternally consistent state after processing each said file system change.

32
33 40. A system as in claim 26, wherein
34 said mass storage element includes a file storage system;

1 each said file server is disposed for leaving said file storage system in an inter-
2 nally consistent state after processing file system changes;
3 said internally consistent state is associated with a set of completed file system
4 changes;
5 said set of completed file system changes is identifiable by each said file server.
6

7 41. A file server system as in claim 26, wherein
8 said mass storage element includes a primary mass storage element and a mirror
9 mass storage element; and
10 said first file server processes said file system change for both said primary mass
11 storage element and said mirror mass storage element.
12

13 42. A method of operating a file server system, said method including steps
14 for
15 responding to an incoming file server request by copying a descriptor of a file
16 system change to both a first file server and a second file server;
17 processing said file system change at said first file server while maintaining said
18 descriptor copy at said second file server; and
19 performing, at said second file server, a file system change in response to a cop-
20 ied descriptor and a service interruption by said first file server.
21

22 43. A method as in claim 42, including steps for associating a first file server
23 and a second file server with a mass storage element.
24

25 44. A method as in claim 42, including steps for delaying output by at least
26 one said file server to said mass storage system without delaying a response to file system
27 changes.
28

29 45. A method as in claim 42, wherein a first said file server is disposed for
30 processing said file system changes atomically, whereby a second said file server can on failover
31 process exactly those file system changes not already processed by said first file server.
32

33 46. A method as in claim 42, wherein a first said file server is disposed to re-
34 spond identically to service interruptions for itself and for a second said file server.
35

1 47. A method as in claim 42, wherein at least one said file server responds to
2 a file system change before committing a result of said file system change to mass storage.

3
4 48. A method as in claim 42, wherein
5 each said file server includes a file server change memory;
6 each one of said file servers is coupled to at least a portion of said file server
7 change memory using local memory access; and
8 each one of said file servers is coupled to at least a portion of said file server re-
9 quest memory using remote memory access.

10
11 49. A method as in claim 42, wherein said file server change memory in-
12 cludes a disk block.

13
14 50. A method as in claim 42, wherein said file server change memory in-
15 cludes a file server request.

16
17 51. A method as in claim 42, wherein said file server change memory is dis-
18 posed to delay output to said mass storage element without delaying a response to file server re-
19 quests.

20
21 52. A method as in claim 42, wherein said mass storage element includes a
22 file storage system and each said file server is disposed for leaving said file storage system in an
23 internally consistent state after processing each said file system change.

24
25 53. A method as in claim 42, wherein said steps for performing a file system
26 change in response to a copied descriptor are also operative in response to a service interruption
27 by said second file server.

28
29 54. A method as in claim 42, wherein said steps for processing includes steps
30 for processing said file system change at both a primary mass storage element and a mirror mass
31 storage element.

32
33 55. A method of operating a file server system, said method including steps
34 for

1 receiving a file server request at one of a plurality of file servers and in response
2 thereto copying a descriptor of a file system change into a server change memory;
3 processing said file system change for both a first set of mass storage devices
4 coupled to a first one said file server and for at least one said mass storage device in a second set
5 of mass storage devices coupled to a second one said file server.

6
7 56. A method as in claim 56, wherein said descriptor includes a file server
8 request.

9
10 57. A method as in claim 56, wherein said server change memory includes a
11 disk block.

12
13 58. A method as in claim 56, wherein said server change memory includes a
14 file server request.

15
16 59. A method as in claim 56, wherein said server change memory includes a
17 first portion disposed at said first file server and a second portion disposed at said second file
18 server.

19
20 60. A method as in claim 56, wherein said server change memory includes a
21 first portion disposed at said first file server and a second portion disposed at said second file
22 server; and wherein said steps for copying include steps for copying said descriptor into both
23 said first portion and said second portion.

24
25 61. A method as in claim 56, wherein said server change memory includes a
26 first portion disposed at said first file server and a second portion disposed at said second file
27 server; and said steps for copying include steps for copying said descriptor into both said first
28 portion and said second portion by either of said first file server or said second file server.

29
30 62. A method as in claim 56, wherein said server change memory is disposed
31 to delay output to said mass storage element without delaying a response to file server requests.

32
33 63. A method as in claim 56, wherein

1 said steps for receiving include receiving a file server request at either said first
2 file server or said second file server, and said steps for copying said descriptor include copying
3 by either said first file server or said second file server; and including steps for
4 processing said file system change for both said second set of mass storage de-
5 vices and for at least one said mass storage device in said first set.

6
7 64. A method of operating a file server system, said method including steps
8 for
9 receiving a file server request at one of a plurality of file servers and in response
10 thereto copying a descriptor of a file system change into a file server change memory; and
11 responding to a service interruption by performing a file system change in re-
12 sponse to a descriptor in said file server change memory.

13
14 65. A method as in claim 65, including steps for associating a plurality of file
15 servers with at least one mass storage element and at least one file server change memory.

16
17 66. A method as in claim 65, including steps for delaying output to a mass
18 storage element without delaying a response to file server requests.

19
20 67. A method as in claim 65, including steps for leaving a file storage system
21 on said mass storage element in an internally consistent state after processing each said file sys-
22 tem change.

23
24 68. A method as in claim 65, including steps for
25 leaving a file storage system on said mass storage element in an internally con-
26 sistent state after processing file system changes;
27 associating said internally consistent state with a set of completed file system
28 changes; and
29 identifying said set of completed file system changes by at least one said file
30 server.

31
32 69. A method as in claim 65, including steps for performing said received file
33 server request at both a primary mass storage element and a mirror mass storage element.

34
35 70. A method as in claim 65, including steps for

1 processing said file system changes atomically at a first said file server; and
2 on failover processing exactly those file system changes not already processed by
3 said first file server.

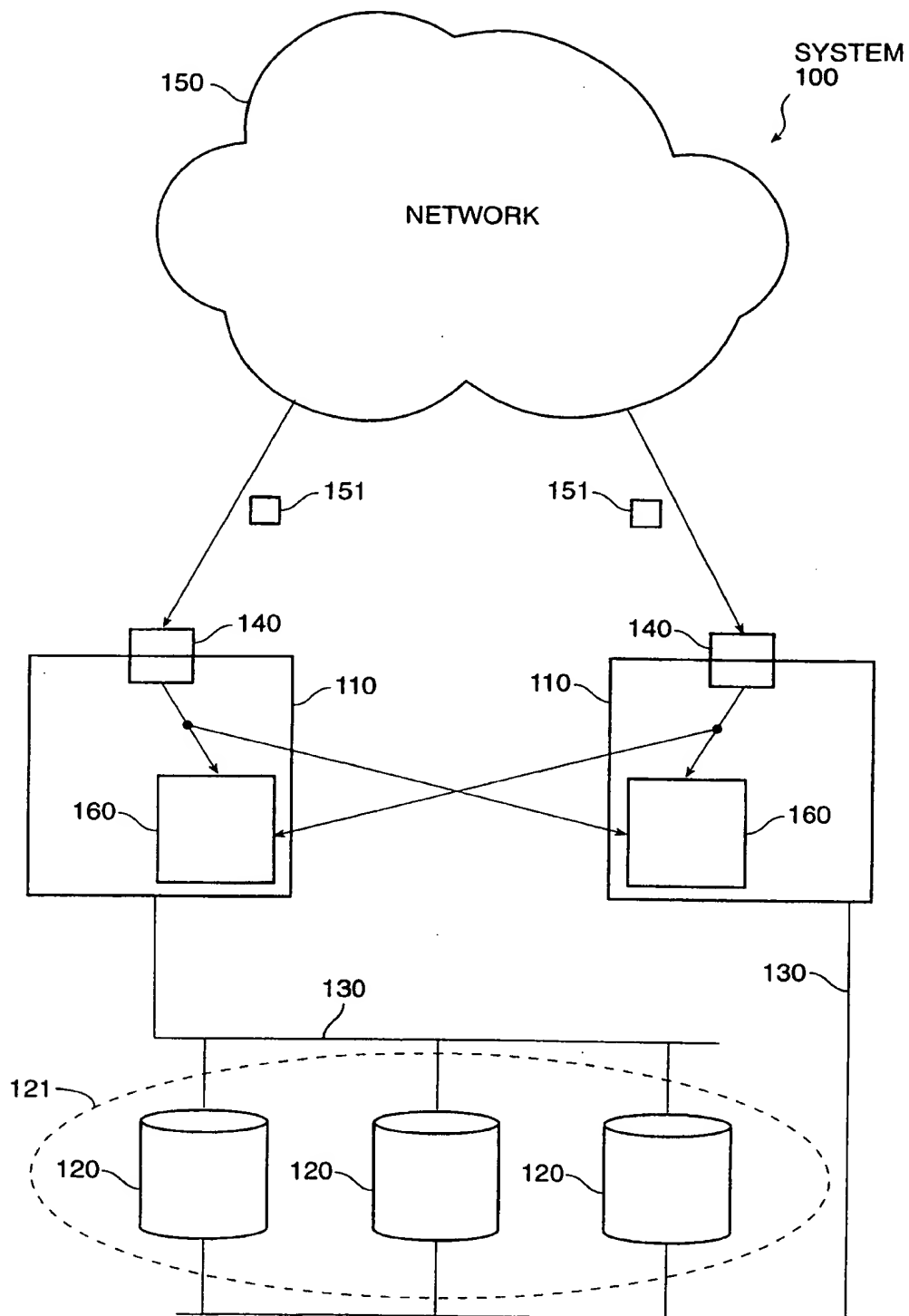
4
5 71. A method as in claim 65, including steps for responding identically at a
6 first said file server to service interruptions for itself and for a second said file server.

7
8 72. A method as in claim 65, wherein said file server change memory in-
9 cludes a disk block.

10
11 73. A method as in claim 65, wherein said file server change memory in-
12 cludes a file server request.

13
14 74. A method as in claim 65, wherein said file server change memory is dis-
15 posed to delay output to said mass storage element without delaying a response to file server re-
16 quests.

17
18 75. A method as in claim 65, including steps for responding to a file system
19 change before committing a result of said file system change to mass storage at one said file
20 server.

**FIG. 1**

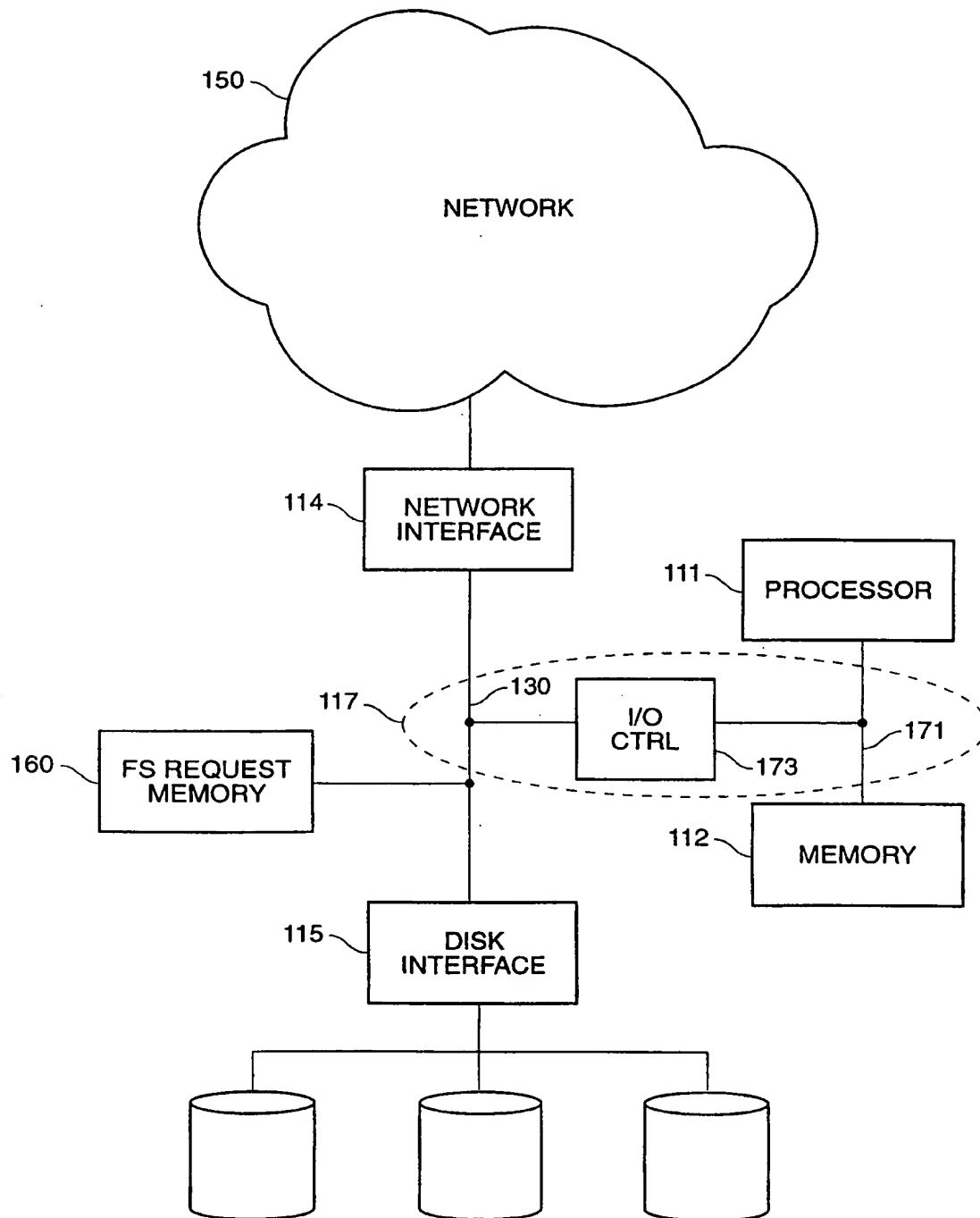


FIG. 2

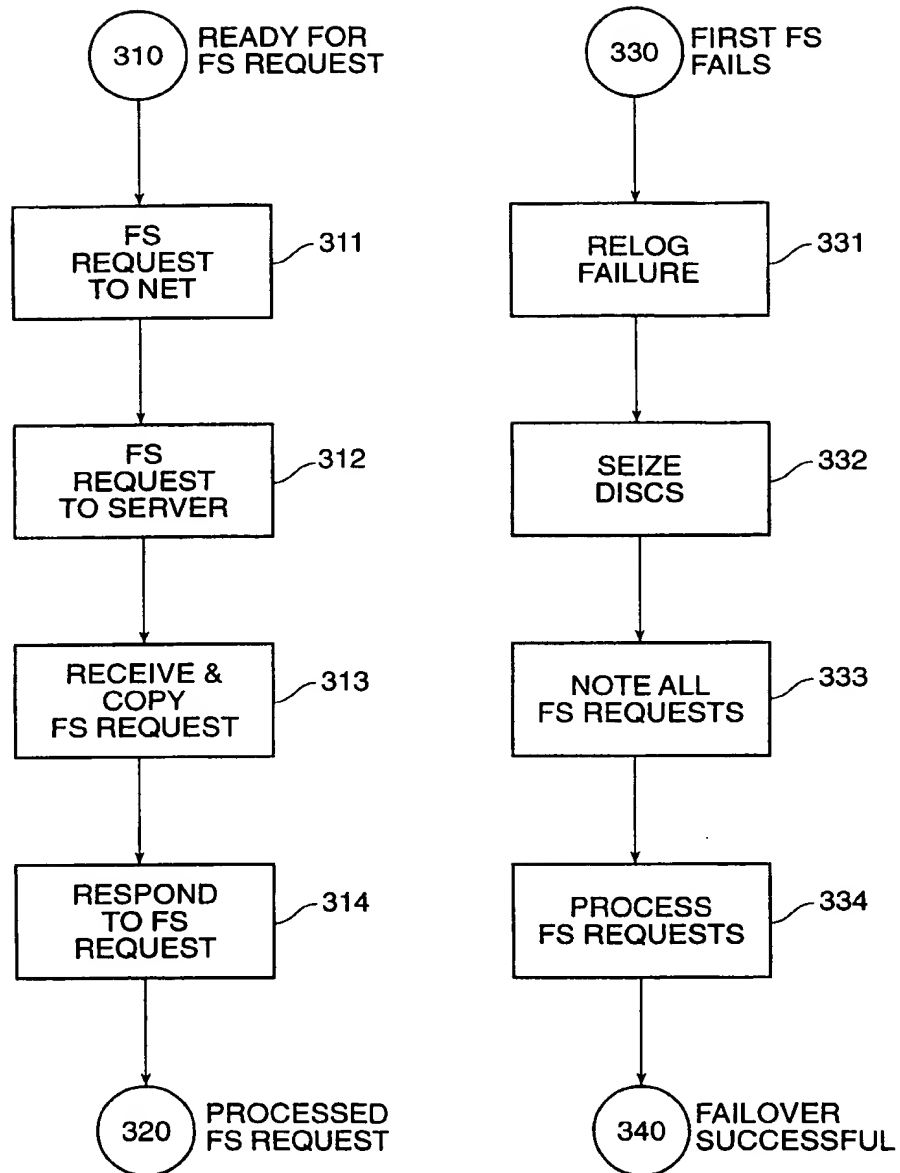
METHOD
300

FIG. 3

**This Page is Inserted by IFW Indexing and Scanning
Operations and is not part of the Official Record**

BEST AVAILABLE IMAGES

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images include but are not limited to the items checked:

- ☐ BLACK BORDERS
- ☐ IMAGE CUT OFF AT TOP, BOTTOM OR SIDES
- ☒ FADED TEXT OR DRAWING
- ☒ BLURRED OR ILLEGIBLE TEXT OR DRAWING
- ☐ SKEWED/SLANTED IMAGES
- ☐ COLOR OR BLACK AND WHITE PHOTOGRAPHS
- ☐ GRAY SCALE DOCUMENTS
- ☒ LINES OR MARKS ON ORIGINAL DOCUMENT
- ☒ REFERENCE(S) OR EXHIBIT(S) SUBMITTED ARE POOR QUALITY
- ☐ OTHER: _____

IMAGES ARE BEST AVAILABLE COPY.

As rescanning these documents will not correct the image problems checked, please do not report these problems to the IFW Image Problem Mailbox.